

## IAT 2 - Natural Language Processing (NLP) - Question Paper

USN



### Internal Assessment Test 2 – May 2025

Sub:	Natural Language Processing					Sub Code:	BAD613B	Branch:	AI&DS		
Date:	23/05/2025	Duration:	90 mins	Max Marks:	50	Sem / Sec:	VI			OBE	
<u>Answer any FIVE FULL Questions</u>									MARKS	CO	RBT
1a	Discuss the applications of Text Classification							[5]	CO2	L2	
1b	Explain multinomial Naive Bayes Classifier and its mathematical principles							[5]	CO2	L3	
2a	Training Data			Test Data				[5]	CO2	L2	
	Message		Label	Message		Label					
	Just plain boring		-	Predictable with no fun		?					
	Entirely predictable and lacks energy		-	Classify the test message as "+" or "-" using Naive Bayes.							
	No surprises and very few laughs		-								
	Very powerful		+								
	The most fun film of the summer		+								
2b	Demonstrate the application of the Naïve Bayes algorithm as a language model and discuss its evaluation metrics.							[5]	CO2	L3	
3a	Write a short notes on Information Retrieval (IR) models							[5]	CO3	L2	
3b	Explain Boolean model with example							[5]	CO3	L3	
4a	What are lexical resources, and how are WordNet, FrameNet, stemmers, POS taggers, and corpora used?							[5]	CO3	L2	
4b	Summarize about Major issues in Information Retrieval.							[5]	CO4	L3	
5a	Explain machine translation using the encoder-decoder architecture.							[5]	CO4	L2	
5b	Discuss language divergences and linguistic typology.							[5]	CO4	L3	
6a	What is machine translation evaluation? Explain its importance.							[5]	CO4	L2	
6b	Describe the issues of bias and ethics in the context of machine translation.							[5]	CO4	L3	

## NLP Answer Key – IAT 2

Sub:	Natural Language Processing			Sub Code:	BAD613B	Branch:	AI&DS
Date :	24/05/2025	Duration:	90 mins	Max Marks:	50	Sem	VI
1a	<p><b><u>Discuss the applications of Text Classification</u></b></p> <ol style="list-style-type: none"> <li><b>Spam Detection:</b> Classifies emails as <i>spam</i> or <i>not spam</i> based on content and keywords.</li> <li><b>Sentiment Analysis:</b> Determines the <i>emotional tone</i> (positive, negative, neutral) in reviews, social media posts, etc.</li> <li><b>News Categorization:</b> Automatically classifies news articles into categories like <i>sports, politics, technology</i>, etc.</li> <li><b>Language Detection:</b> Identifies the <i>language</i> of a given text (e.g., English, French, Tamil).</li> <li><b>Topic Labeling in Customer Support:</b> Tags customer queries or tickets by topic (e.g., <i>billing, technical issue, feedback</i>) for quicker resolution.</li> </ol>						
1b	<p><b><u>Explain multinomial Naive Bayes Classifier and its mathematical principles</u></b></p> <ol style="list-style-type: none"> <li><b>Naïve Bayes</b> is a classification algorithm based on <b>Bayes' Theorem</b>, used to predict the class of given data.</li> <li>It assumes that all input features are <b>independent</b> of each other — this is the "naïve" assumption.</li> <li><b>Bayes' Theorem:</b></li> </ol> $P(C X) = \frac{P(X C) \cdot P(C)}{P(X)}$ <p>Where <math>P(C X)</math> is the probability of class C given data X.</p> <ol style="list-style-type: none"> <li>The classifier selects the class with the <b>highest probability</b> using:</li> </ol> $\hat{C} = \underset{C_k}{arg\ max} P(C_k) \prod_{i=1}^n P(x_i C_k)$ <ol style="list-style-type: none"> <li>It is simple, fast, and works well for <b>text classification</b> problems like spam detection and sentiment analysis.</li> </ol>						

2a

Training Data		Test Data	
Message	Label	Message	Label
Just plain boring	-	Predictable with no fun	?
Entirely predictable and lacks energy	-	Classify the test message as “+” or “-” using Naive Bayes.	
No surprises and very few laughs	-		
Very powerful	+		
The most fun film of the summer	+		

The prior  $P(c)$  for the two classes is computed via Eq. 4.11 as  $\frac{N_c}{N_{doc}}$ :

$$P(-) = \frac{3}{5} \quad P(+) = \frac{2}{5}$$

The word *with* doesn't occur in the training set, so we drop it completely (as mentioned above, we don't use unknown word models for naive Bayes). The likelihoods from the training set for the remaining three words “predictable”, “no”, and “fun”, are as follows, from Eq. 4.14 (computing the probabilities for the remainder of the words in the training set is left as an exercise for the reader):

$$P(\text{“predictable”}|-) = \frac{1+1}{14+20} \quad P(\text{“predictable”}|+) = \frac{0+1}{9+20}$$

$$P(\text{“no”}|-) = \frac{1+1}{14+20} \quad P(\text{“no”}|+) = \frac{0+1}{9+20}$$

$$P(\text{“fun”}|-) = \frac{0+1}{14+20} \quad P(\text{“fun”}|+) = \frac{1+1}{9+20}$$

For the test sentence  $S = \text{“predictable with no fun”}$ , after removing the word ‘with’ the chosen class, via Eq. 4.9, is therefore computed as follows:

$$P(-)P(S|-) = \frac{3}{5} \times \frac{2 \times 2 \times 1}{34^3} = 6.1 \times 10^{-5}$$

$$P(+)P(S|+) = \frac{2}{5} \times \frac{1 \times 1 \times 2}{29^3} = 3.2 \times 10^{-5}$$

The model thus predicts the class *negative* for the test sentence.

2b

**Demonstrate the application of the Naïve Bayes algorithm as a language model and discuss its evaluation metrics.**

- Language Modeling:** Naïve Bayes can be used to model the probability of words in a language and classify text into specific languages based on word occurrence.
- Text Classification:** It is widely used for tasks like spam detection, sentiment analysis, and language identification using word frequency and class probability.
- Training Process:** The algorithm calculates prior probabilities  $P(C)$  for each language and likelihood  $P(w_i | C)$  for each word  $w_i$  given a class (language).

	<p>4. <b>Evaluation Metrics:</b> Common metrics include:</p> <p><b>Accuracy:</b> Overall correctness of the model.</p> <p><b>Precision:</b> Correct positive predictions out of total predicted positives.</p> <p><b>Recall:</b> Correct positive predictions out of actual positives.</p> <p><b>F1-Score:</b> Harmonic mean of precision and recall.</p> <p>5. <b>Advantage:</b> Naïve Bayes is simple, fast, and effective for large-scale language-based classification problems with high-dimensional data.</p>
3a	<p><b><u>Write a short notes on Information Retrieval (IR) models</u></b></p> <p>1. <b>Boolean Model (Classical):</b> Uses logical operators (AND, OR, NOT) for exact match retrieval. Simple but does not rank documents.</p> <p>2. <b>Vector Space Model (Classical):</b> Represents documents and queries as vectors and uses <b>cosine similarity</b> to rank based on relevance.</p> <p>3. <b>Probabilistic Model / BM25 (Classical):</b> Estimates the probability of document relevance to a query. BM25 is a popular ranking function based on this model.</p> <p>4. <b>Fuzzy and Extended Boolean Models (Non-Classical):</b> Allow partial matches using weighted terms or fuzzy logic, improving flexibility over standard Boolean retrieval.</p> <p>5. <b>Language Models (Non-Classical):</b> Treat each document as a probability distribution and rank by the likelihood that it generates the query.</p> <p>6. <b>Latent Semantic Indexing (LSI) (Alternative):</b> Uses mathematical techniques (like SVD) to identify hidden relationships between terms and documents.</p> <p>7. <b>Neural IR Models (Alternative):</b> Use deep learning (e.g., BERT) for semantic matching and contextual understanding, offering state-of-the-art performance.</p>

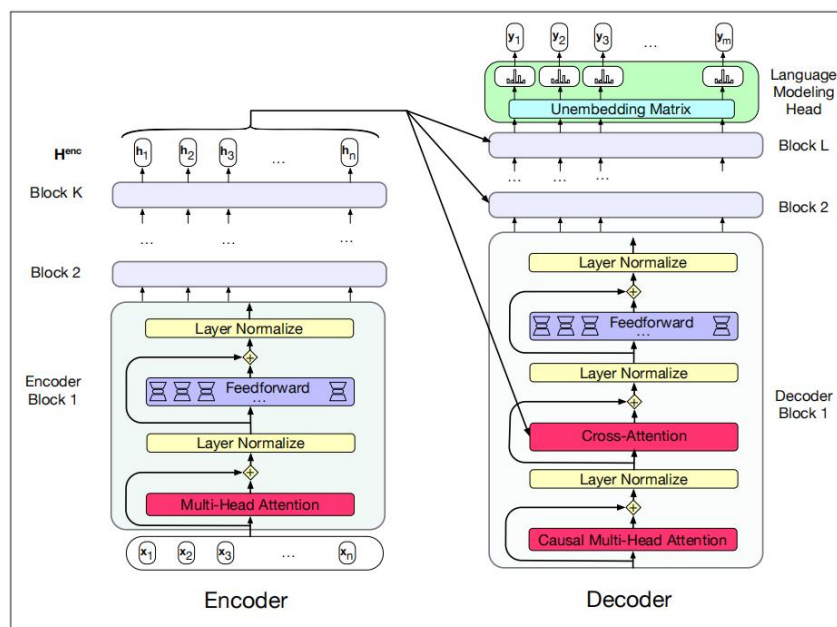
3b	<p><b><u>Explain Boolean model with example</u></b></p> <p><b>1. Definition:</b> The Boolean model is a classical Information Retrieval model that represents documents and queries as sets of keywords. It uses <b>Boolean logic operators</b> like <b>AND, OR, NOT</b> to retrieve documents that match the query exactly.</p> <p><b>2. Working Principle:</b> Each document is either <b>relevant (1)</b> or <b>not relevant (0)</b> — no ranking is involved. Queries are formed using Boolean expressions to combine terms.</p> <p><b>3. Operators:</b></p> <p style="padding-left: 40px;"><b>AND:</b> Retrieves documents containing <i>all</i> specified terms.</p> <p style="padding-left: 40px;"><b>OR:</b> Retrieves documents containing <i>any</i> of the terms.</p> <p style="padding-left: 40px;"><b>NOT:</b> Excludes documents containing a specific term.</p> <p><b>4. Example:</b> Suppose we have a document set:</p> <p style="padding-left: 40px;">D1: "apple banana mango"</p> <p style="padding-left: 40px;">D2: "banana orange"</p> <p style="padding-left: 40px;">D3: "apple mango"</p> <p style="padding-left: 40px;">Query: apple AND mango</p> <p style="padding-left: 40px;">Result: D1 and D3 (both contain "apple" and "mango")</p> <p><b>5. Advantages and Limitations:</b></p> <p style="padding-left: 40px;"><b>Pros:</b> Simple to understand and implement.</p> <p style="padding-left: 40px;"><b>Cons:</b> No ranking of documents; can't handle partial relevance or term weighting.</p>
4a	<p><b><u>What are lexical resources, and how are WordNet, FrameNet, stemmers, POS taggers, and corpora used?</u></b></p> <p><b>1. Lexical Resources:</b> These are databases or tools that provide information about words, their meanings, relationships, usage, and grammatical properties. They are essential in Natural Language Processing (NLP) tasks like parsing, tagging, and semantic analysis.</p> <p><b>2. WordNet:</b> A lexical database of English where words are grouped into sets of synonyms (synsets) with definitions and semantic relations (e.g., synonyms, antonyms, hypernyms). Used in word sense disambiguation and semantic similarity tasks.</p> <p><b>3. FrameNet:</b> A resource based on frame semantics. It documents how words evoke specific</p>

	<p>situations (frames) and the roles (participants) in them. Used in semantic role labeling and deep NLP understanding.</p> <p>4. Stemmers: Tools that reduce words to their root form by removing suffixes (e.g., “running” → “run”). Common stemmers: Porter, Snowball. Used in information retrieval and text normalization.</p> <p>5. POS Taggers (Part-of-Speech Taggers): Tools that label words with their grammatical category (noun, verb, adjective, etc.). Useful in syntactic parsing, chunking, and named entity recognition.</p> <p>6. Corpora: Large and structured collections of text (e.g., Brown Corpus, Penn Treebank) used to train, test, and evaluate NLP models.</p>
4b	<p><b><u>Summarize about Major issues in Information Retrieval.</u></b></p> <ul style="list-style-type: none"> <li>✓ <b>Query Understanding:</b> Users often submit short or ambiguous queries. Interpreting the actual intent behind a query is a major challenge in IR.</li> <li>✓ <b>Relevance Determination:</b> Determining what makes a document <i>relevant</i> to a query varies between users. IR systems struggle with modeling subjective relevance accurately.</li> <li>✓ <b>Ranking of Results:</b> It is difficult to design effective ranking algorithms that consistently place the most relevant documents at the top.</li> <li>✓ <b>Handling Synonyms and Polysemy:</b></li> <li>✓ <b>Synonyms:</b> Different words with the same meaning (e.g., “car” and “automobile”).</li> <li>✓ <b>Polysemy:</b> Same word with different meanings (e.g., “bank” – riverbank or financial institution). Both cause problems in matching queries to documents.</li> <li>✓ <b>Scalability and Performance:</b> With massive volumes of data on the web, IR systems must be efficient, fast, and scalable for real-time performance.</li> </ul>
5a	<p><b><u>Explain machine translation using the encoder-decoder architecture.</u></b></p> <p>Machine Translation is the automatic conversion of text from one language to another. The <b>encoder-decoder architecture</b> is a popular deep learning approach used for this task.</p> <ol style="list-style-type: none"> <li>1. <b>Encoder:</b> The encoder processes the input sentence (source language) word by word using models like RNN, LSTM, or GRU. It converts the entire sentence into a fixed-length context vector (a summary of the input).</li> <li>2. <b>Context Vector:</b> This vector captures the meaning of the whole sentence and is passed to the decoder.</li> <li>3. <b>Decoder:</b> The decoder takes the context vector and generates the output sentence (target language) one word at a time. It predicts the next word based on the context and the previously generated words.</li> <li>4. <b>Training:</b> The model is trained using parallel corpora (pairs of sentences in</li> </ol>

source and target languages) to minimize translation errors.

5. **Limitations and Improvements:** A basic encoder-decoder struggles with long sentences. To improve this, **attention mechanisms** are used, allowing the decoder to focus on relevant parts of the input sentence at each step.

This architecture forms the foundation for many modern neural machine translation systems.



The transformer block for the encoder and the decoder.

The transformer-based encoder-decoder architecture for machine translation consists of two main parts: the **encoder**, which processes the source sentence into contextual embeddings ( $H^{enc}$ ), and the **decoder**, which generates the target sentence one word at a time. The decoder includes an extra **cross-attention layer** that allows it to attend to all encoder outputs, using them as keys and values while generating each word. This attention helps align source and target tokens effectively. Training is done using **teacher forcing** and **cross-entropy loss**, where the model is trained to predict the next word using the correct previous target token. Decoding often uses **beam search** for better translation quality.

5b

**Discuss language divergences and linguistic typology.**

- **Language Divergences** refer to the structural and semantic differences between languages that can make machine translation challenging. These include:
  - ✓ **Lexical divergence:** One language may have a single word, while another may need a phrase (e.g., “uncle” in English vs. separate words for maternal/paternal uncle in Hindi).
  - ✓ **Syntactic divergence:** Differences in word order, like Subject-Verb-

	<p>Object (SVO) in English vs. Subject-Object-Verb (SOV) in Japanese.</p> <ul style="list-style-type: none"> <li>✓ <b>Morphological divergence:</b> Some languages use rich inflection (e.g., Turkish), while others use word order or prepositions.</li> </ul> <p>➤ <b>Linguistic Typology</b> is the classification of languages based on their common structural features. It helps in understanding:</p> <ul style="list-style-type: none"> <li>✓ Word order types (SVO, SOV, etc.)</li> <li>✓ Morphological types (isolating, agglutinative, fusional)</li> <li>✓ Syntactic and grammatical patterns</li> </ul>
6a	<p><b><u>What is machine translation evaluation? Explain its importance.</u></b></p> <p><b>Machine Translation (MT) Evaluation</b> is the process of assessing the quality of automatically translated text. It determines how well the translated output preserves the meaning, fluency, and grammatical correctness of the source text.</p> <p>There are two main types:</p> <ul style="list-style-type: none"> <li>✓ <b>Automatic Evaluation:</b> Uses metrics like BLEU, METEOR, and ROUGE to compare machine output with human reference translations.</li> <li>✓ <b>Human Evaluation:</b> Involves human judges rating translations based on adequacy (meaning preservation) and fluency (naturalness of language).</li> </ul> <p><b>Importance:</b></p> <ul style="list-style-type: none"> <li>✓ Ensures the <b>accuracy</b> and <b>reliability</b> of translations.</li> <li>✓ Helps in <b>comparing different MT systems</b>.</li> <li>✓ Guides <b>model improvement</b> by identifying errors.</li> </ul> <p>Essential for real-world applications like healthcare, legal, and international communication where precision is critical.</p>
6b	<p><b><u>Describe the issues of bias and ethics in the context of machine translation.</u></b></p> <p>Machine Translation (MT) systems can unintentionally reflect and amplify <b>biases</b> present in their training data. These biases may be related to <b>gender</b>, <b>culture</b>, or <b>social norms</b>. For example, translating gender-neutral words from one language to another may result in gender-stereotyped outputs (e.g., translating a neutral profession in a source language to “he is a doctor” or “she is a nurse”).</p> <p><b>Ethical issues</b> include:</p> <ul style="list-style-type: none"> <li>✓ <b>Fairness:</b> MT should treat all groups equally and avoid reinforcing stereotypes.</li> <li>✓ <b>Transparency:</b> Users should be informed when translations are</li> </ul>



	<p>machine-generated.</p> <ul style="list-style-type: none"><li>✓ <b>Accountability:</b> Developers must ensure responsible use, especially in sensitive domains like healthcare or law.</li><li>✓ <b>Data Privacy:</b> MT systems should not leak or misuse user data.</li></ul> <p>Addressing bias and ethics is crucial to ensure that MT tools are <b>inclusive, fair, and trustworthy</b> in global communication.</p>
--	---